

# **DYNAMIC PROGRAMMING FOR ACTION IN THE ENVIRONMENT**

Cognitive Dynamic Systems, *Simon Haykin*

Presented by

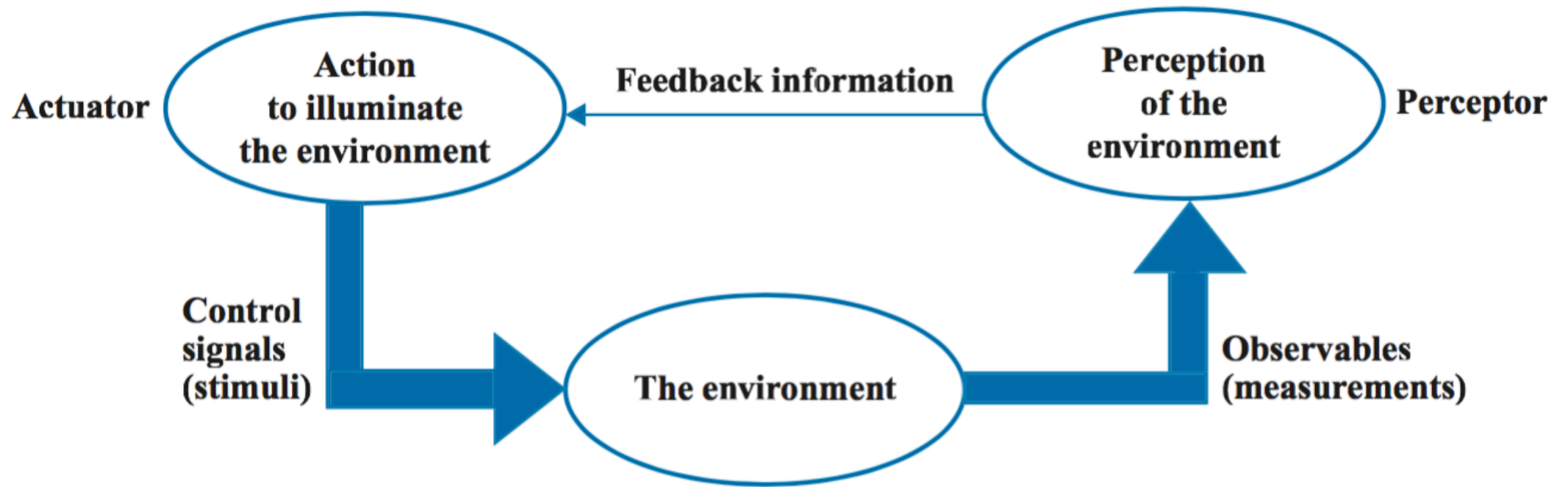
Joo Yeon Kim & Yunseong Lee

# **THE BIG PICTURE**

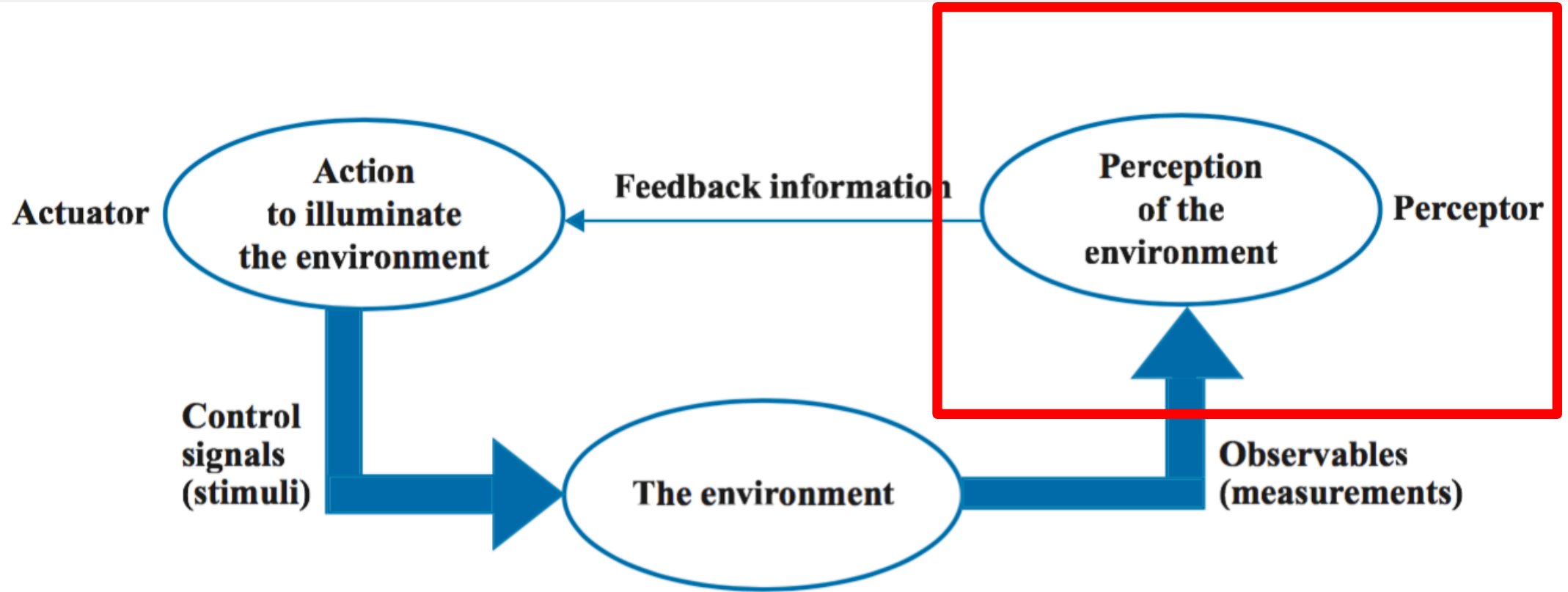
# COGNITIVE DYNAMIC SYSTEMS

Build up rules of behavior over time through learning from continuous experiential interactions with the environment, and thereby deal with environmental uncertainties.

# THE PERCEPTION-ACTION CYCLE



# THE PERCEPTION-ACTION CYCLE



# THE PERCEPTION-ACTION CYCLE

Ac

## Power Spectrum:

Power Spectrum Estimation  
(Cognitive Radio)

## Bayesian Filtering:

State Estimation  
(Cognitive Radar)

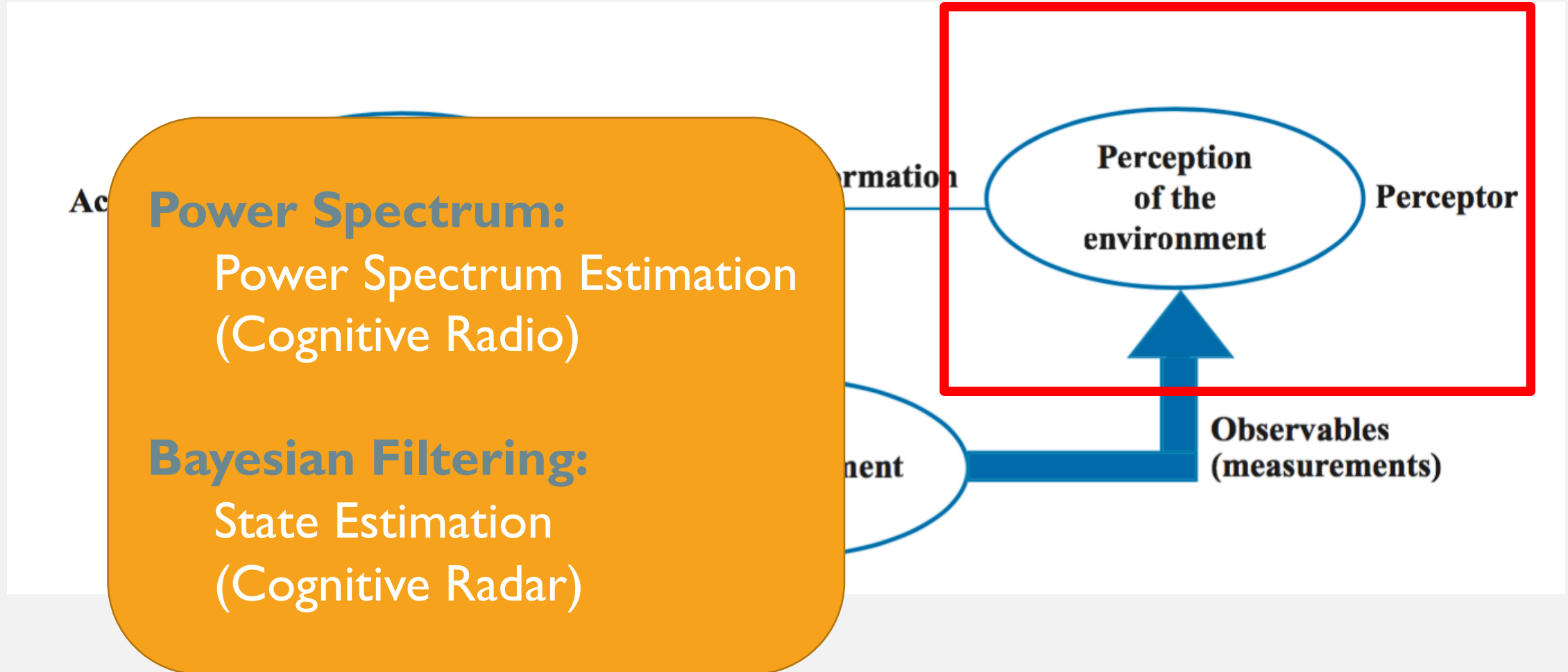
Information

Environment

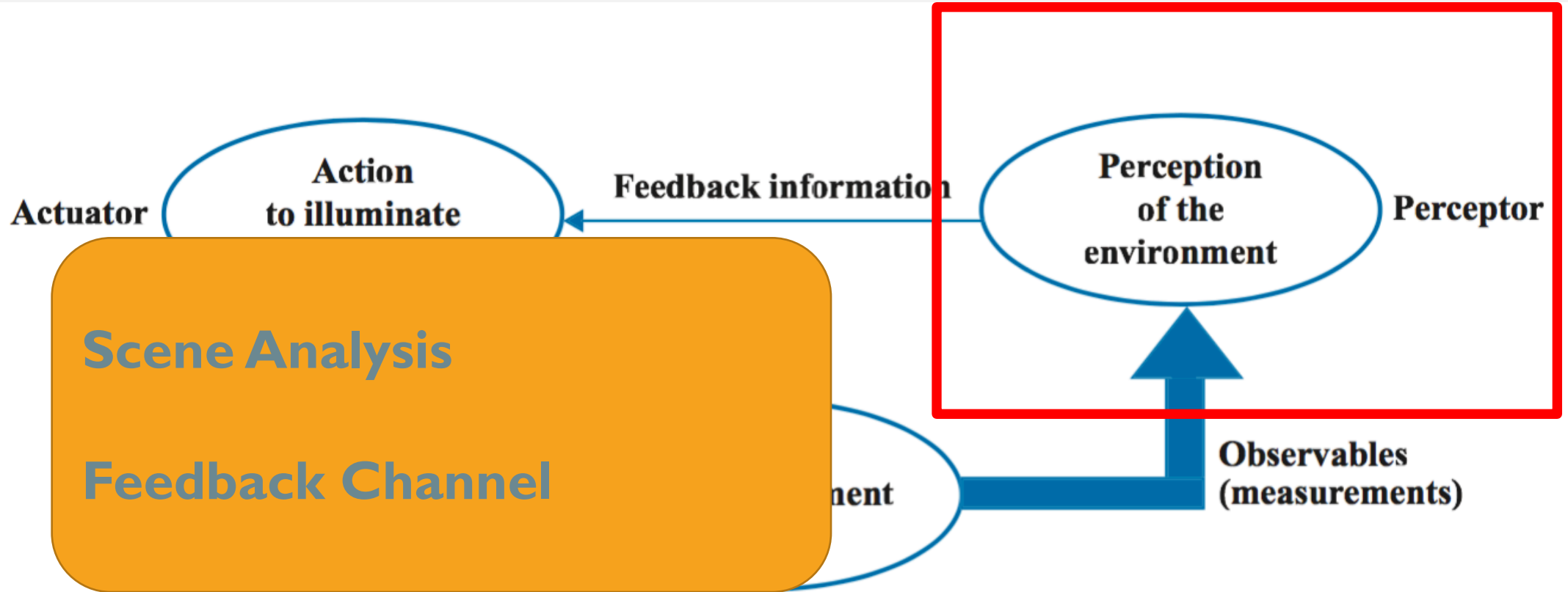
Perception  
of the  
environment

Perceptor

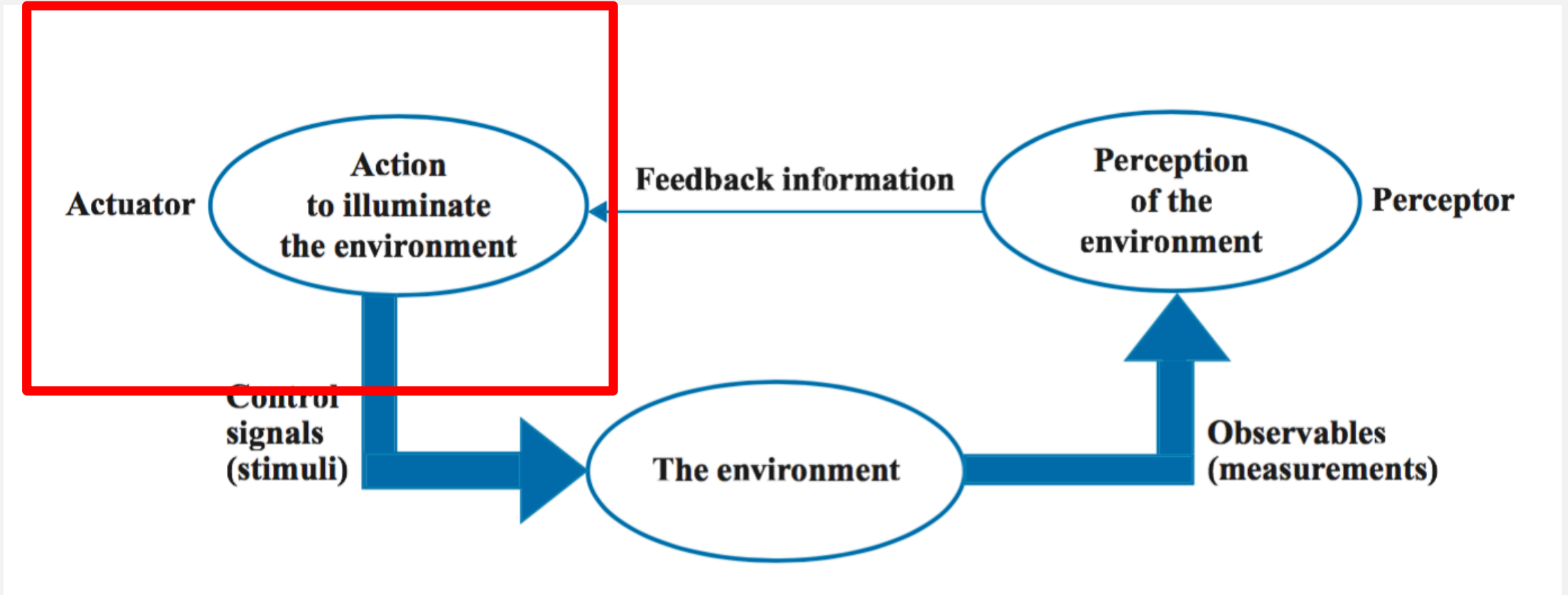
Observables  
(measurements)



# THE PERCEPTION-ACTION CYCLE



# THE PERCEPTION-ACTION CYCLE





# **DYNAMIC PROGRAMMING**

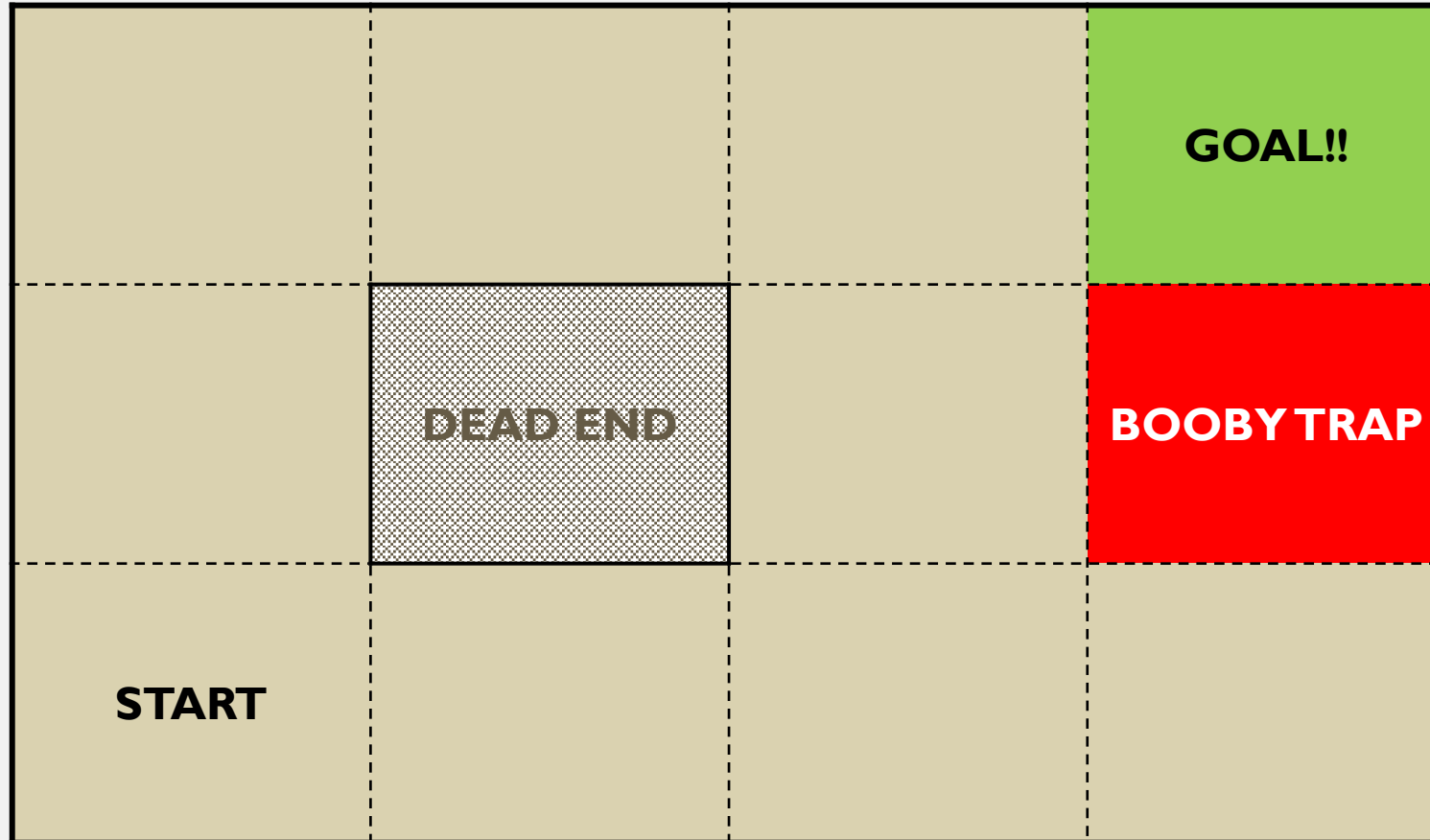
# WHAT IS DYNAMIC PROGRAMMING?

- *Dynamic programming* is a technique that deals with situations where decisions are made in stages (i.e. different time steps), with the outcome of each decision being *predictable* to some extent before the next decision is made.
- A key aspect of such situations is that decisions cannot be made in isolation. Rather, the desire for a low cost at the present is balanced against the undesirability of a high cost in the future.

# DYNAMIC PROGRAMMING

- Markov Decision Processes
- Bellman's Optimality Criterion
- Policy Iteration
- Value Iteration

# MARKOV'S DECISION PROCESS



You can move UP, DOWN, LEFT, RIGHT

# DETERMINISTIC VS. STOCHASTIC

## Deterministic

- Every state can be uniquely determined with model parameters and previous states.
- Always perform identically for a given set of initial conditions.

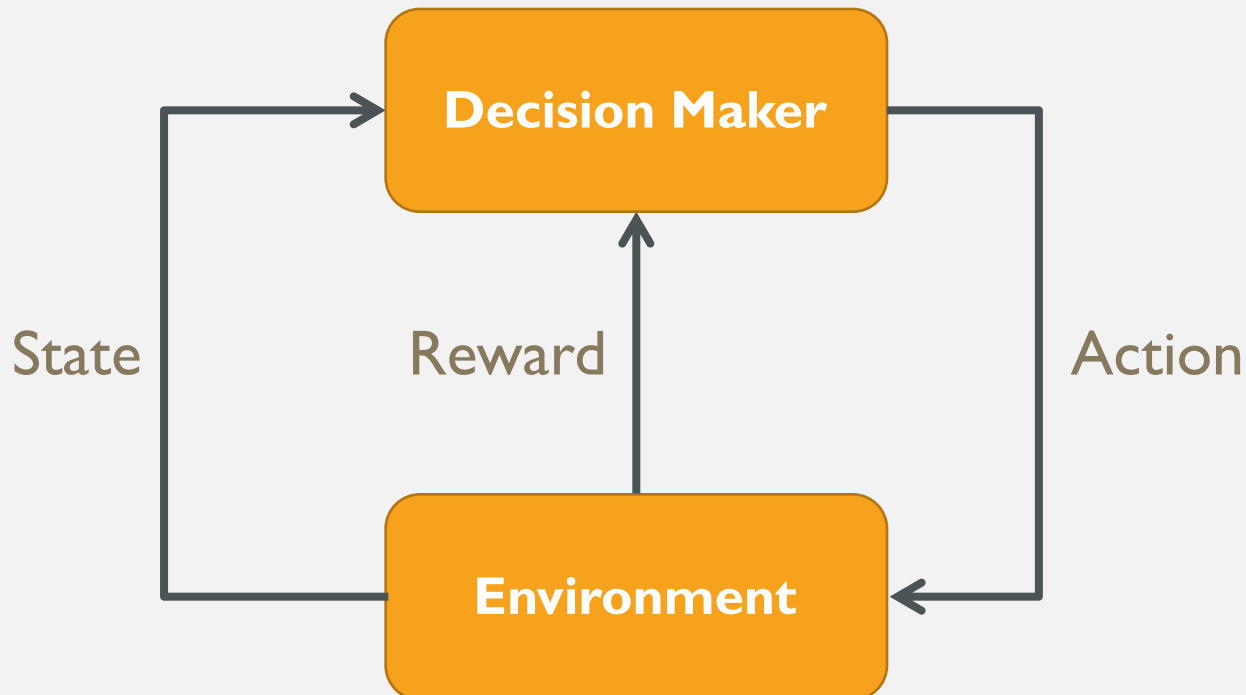
## Stochastic

- States determined by probability distributions.
- Such probability distributions include a notion of randomness.

# MARKOV'S DECISION PROCESSES

## The Actuator's Decision Making Process

*Discrete-time & Stochastic*



- Set of States:  $\mathcal{S} = \{s\}$
- Set of Actions:  $\mathcal{A} = \{a\}, a: \mathcal{S} \rightarrow \mathcal{S}$
- Reward Function:  $R(s)$
- Transition Model:  
 $T(s, a, s')$  or  $P(s' | s, a)$
- Policy:  $\pi(s) \rightarrow a$

# MARKOV'S DECISION PROCESSES

## The Actuator's Decision Making Process

*Discrete-time & Stochastic*

### Markov Property

Decision Maker

$$P(s_{t+1} | a, s_0, \dots, s_t) = P(s_{t+1} | a, s_t)$$

- Set of States:  $S = \{s\}$
- Set of Actions:  $A = \{a\}, a: S \rightarrow S$
- Reward Function:  $R(\cdot)$
- Discount Factor:  $\gamma$
- Initial State:  $s_1$

Environment

State

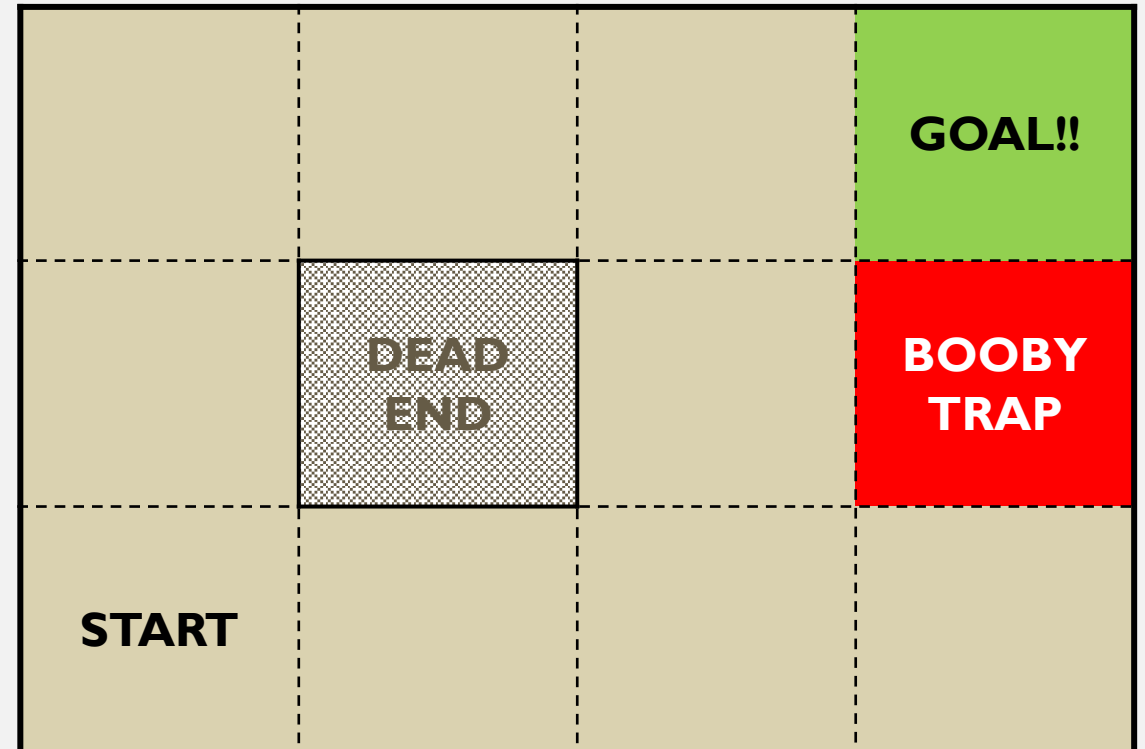
Reward

Action

$S$

# MARKOV'S DECISION PROCESS

- Set of States:  $\mathcal{S} = \{s\}$
- Set of Actions:  $\mathcal{A} = \{a\}, a: \mathcal{S} \rightarrow \mathcal{S}$
- Reward Function:  $R(s)$
- Transition Model:  
 $T(s, a, s')$  or  $P(s' | s, a)$
- Policy:  $\pi(s) \rightarrow a$



You can move UP, DOWN, LEFT, RIGHT

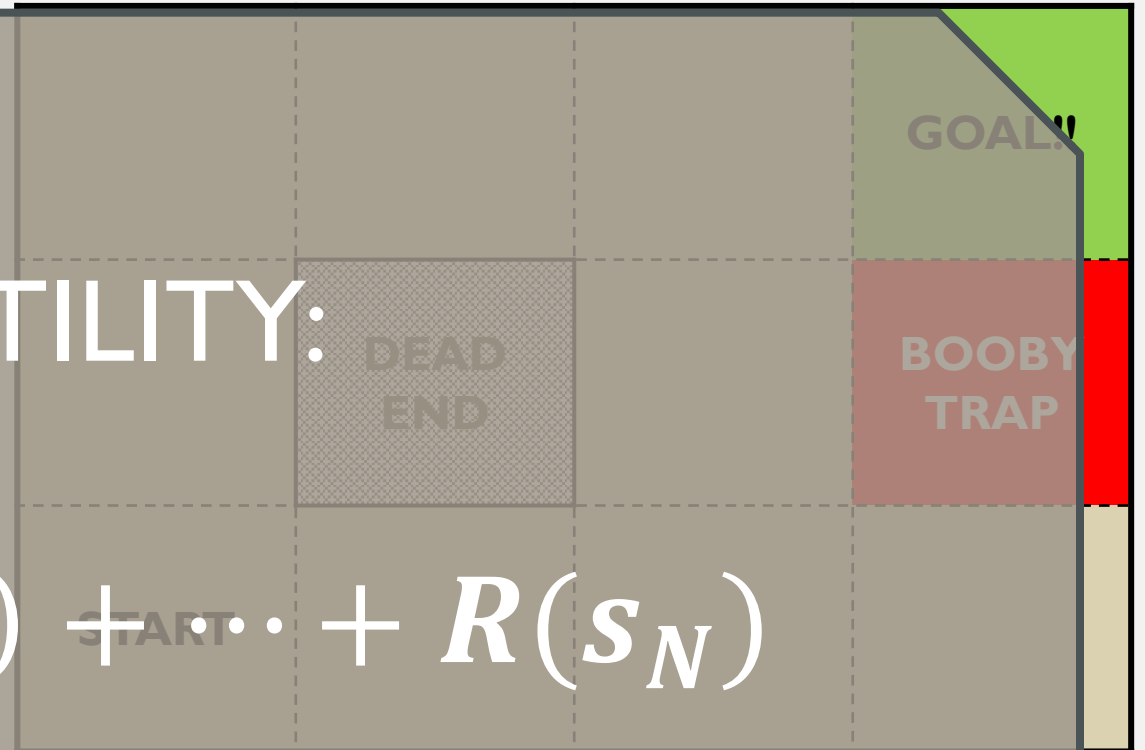


# MARKOV'S DECISION PROCESS

- Set of States:  $\mathcal{S} = \{s\}$
- Set of Actions:  $A = \{a\}, a: \mathcal{S} \rightarrow \mathcal{S}$
- Reward Function:  $R(s)$
- Transition Model:  
 $T(s, a, s')$  or  $P(s' | s, a)$
- Policy:  $\pi(s) \rightarrow a$

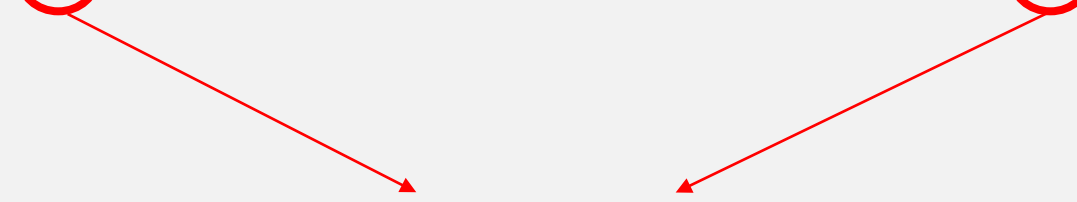
Define UTILITY:

$$U(s_0, \dots, s_N) = R(s_0) + R(s_1) + \dots + R(s_N)$$



You can move UP, DOWN, LEFT, RIGHT

# MARKOV'S DECISION PROCESS

$$U(s_0, \dots, s_N) = R(s_0) + R(s_1) + \dots + R(s_N)$$


But this **N** can increase infinitely!

$$U(s_0, \dots) = R(s_0) + \gamma^1 R(s_1) + \dots + \gamma^N R(s_N) + \dots$$

Discount Factor:  $0 < \gamma < 1$

# MARKOV'S DECISION PROCESS

$$U(\mathbf{s}_0, \dots) = R(\mathbf{s}_0) + \gamma^1 R(\mathbf{s}_1) + \dots + \gamma^N R(\mathbf{s}_N) + \dots$$

We move from  $s_t$  to  $s_{t+1}$  by **action**  $\pi(\mathbf{s}_t)$

## GOAL

Find a **sequence of actions**,  $\pi(\mathbf{s}_t)$ , such that the resulting sequence of states *maximizes* the total discounted reward,  $\mathbf{U}$

# MARKOV'S DECISION PROCESS

$$U(s_0, \dots) = R(s_0) + \gamma^1 R(s_1) + \dots + \gamma^N R(s_N) + \dots$$

## Optimal Policy:

We move from  $s_t$  to  $s_{t+1}$  by action  $\pi(s_t)$

The policy  $\pi^*$  that *maximizes* the expected *utility*  $U$  of the sequence of states generated by  $\pi^*$  for all initial states  $s$

# MARKOV'S DECISION PROCESS

**Question:**

**How do we compute the optimal policy  $\pi^*$ ?**

# BELLMAN'S EQUATION

- If action  $a$  is taken in state  $s$ ,

$$U(s) = R(s, a) + \gamma \sum_{s'} T(s, a, s') U(s')$$

$U(s)$ : Utility in state  $s$

$R(s, a)$ : Reward when taking action  $a$  from state  $s$

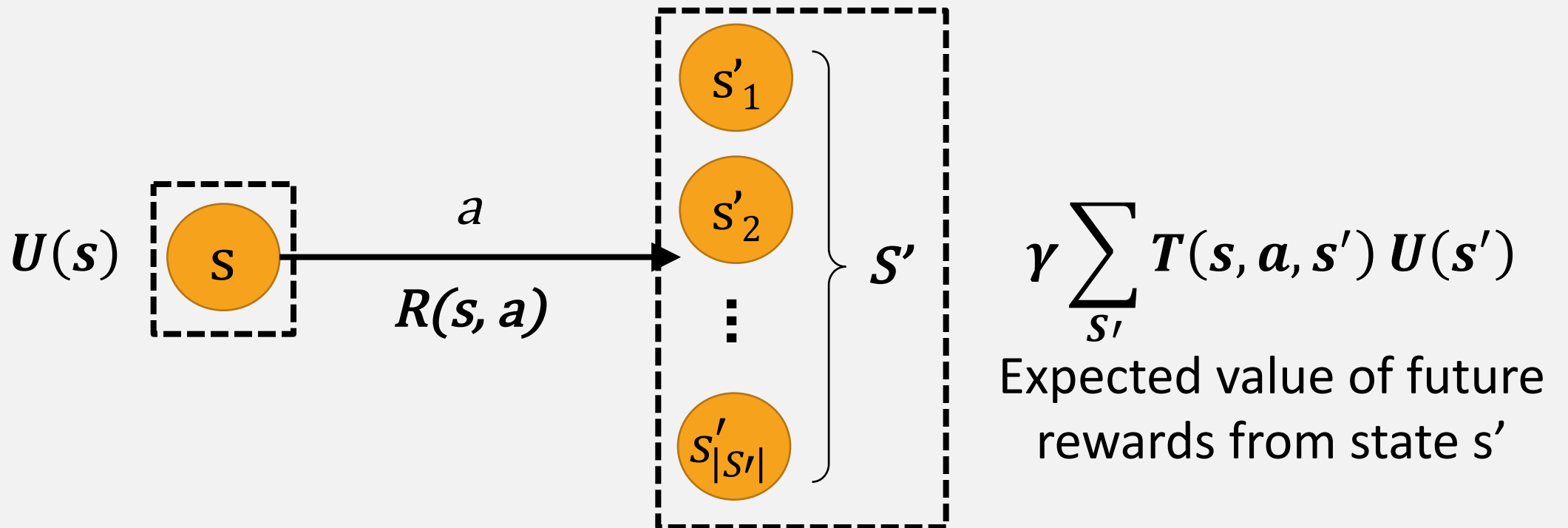
$\gamma$ : Discount factor

$T(s, a, s')$ : Transition probability from  $s$  to  $s'$  by taking  $a$

# BELLMAN'S EQUATION

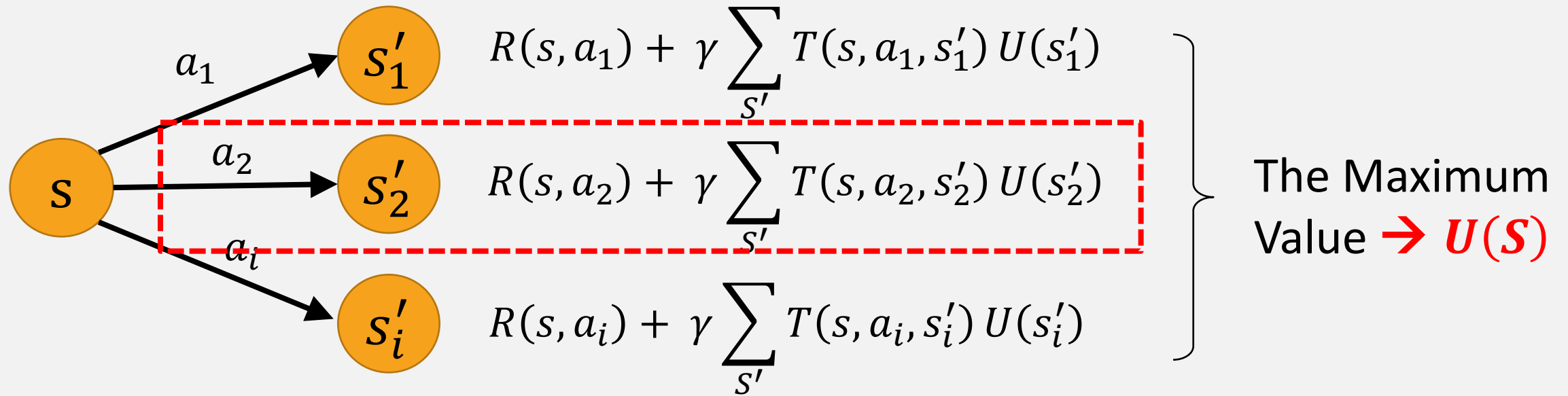
- If action  $a$  is taken in state  $s$ ,

$$U(s) = R(s, a) + \gamma \sum_{s'} T(s, a, s') U(s')$$



# BELLMAN'S OPTIMALITY CRITERION

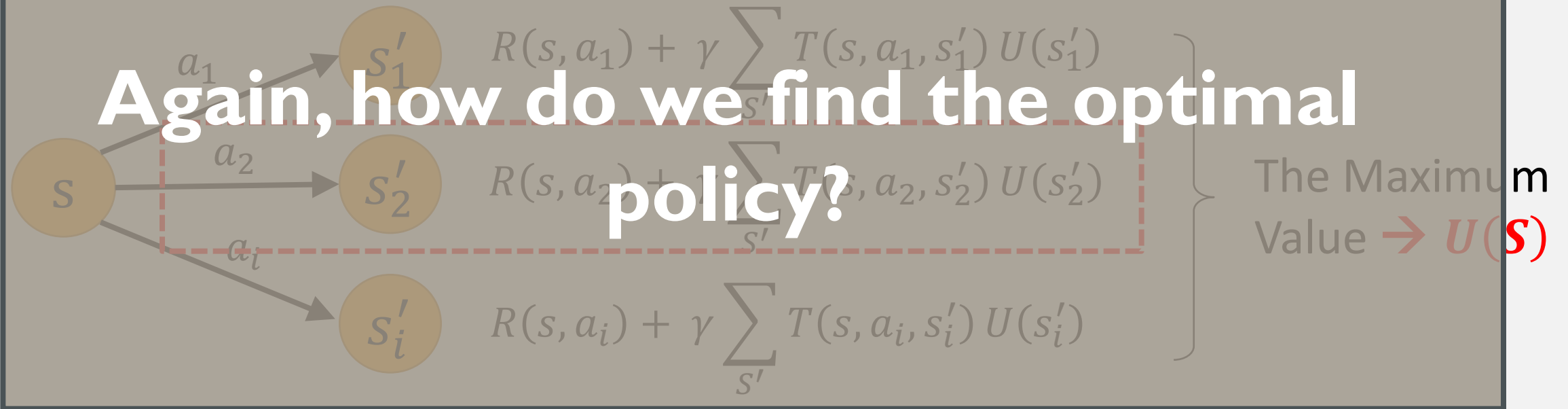
$$U(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') U(s'))$$





# BELLMAN'S OPTIMALITY CRITERION

$$U(S) = \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') U(s'))$$



# POLICY ITERATION

**Initialize:**  $\boldsymbol{\pi}_0 \leftarrow$  guess

**Evaluate:** given  $\boldsymbol{\pi}_t$ ,

calculate  $U_t(\mathbf{s}) = R(\mathbf{s}) + \gamma \sum_{\mathbf{s}'} T(\mathbf{s}, \boldsymbol{\pi}_t(\mathbf{s}), \mathbf{s}') U_t(\mathbf{s}')$  (Bellman's eq)

**Improve:**  $\boldsymbol{\pi}_{t+1} = \operatorname{argmax}_a \sum_{\mathbf{s}'} T(\mathbf{s}, \boldsymbol{\pi}_t(\mathbf{s}), \mathbf{s}') U_t(\mathbf{s}')$

# POLICY ITERATION

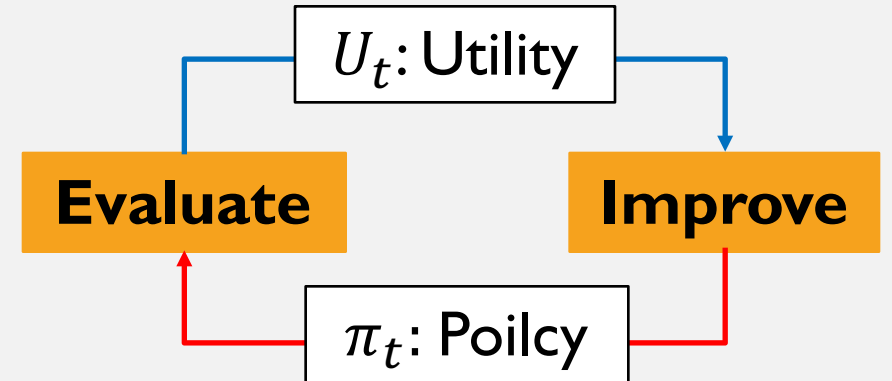
Initialize:  $\pi_0 \leftarrow$  guess

Evaluate: given  $\pi_t$ ,

calculate  $U_t(s) = R(s) + \gamma \sum_{s'} T(s, \pi_t(s), s') U_t(s')$  (Bellman's eq)

Improve:  $\pi_{t+1} = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U_t(s')$

Repeat until the Utility converges  $\rightarrow U_t(s) \approx U_{t+1}(s)$



# VALUE ITERATION

**Initialize:** Start with arbitrary utilities

**Update** utilities based on neighbors

$$\hat{U}_{t+1}(\mathbf{s}) = R(s) + \gamma \max_a \sum_{s'} T(s, \pi_t(s), s') \hat{U}_t(\mathbf{s}') \quad (\text{Bellman's eq})$$

# VALUE ITERATION

**Initialize:** Start with arbitrary utilities

**Update** utilities based on neighbors

$$\hat{U}_{t+1}(\mathbf{s}) = R(\mathbf{s}) + \gamma \max_a \sum_{s'} T(s, \pi_t(s), s') \hat{U}_t(\mathbf{s}') \quad (\text{Bellman's eq})$$

→ Repeat until converges.

The quality of Utility gets better as more *truth* gets involved

# CHECKPOINT

- **Markov Decision Process**

Decision making process in a cognitive dynamic system

- **Bellman's Optimality Criterion**

Defines the optimization problem the decision-maker must solve

- **Policy Iteration**

Finds the optimal policy by iteratively estimating the policy

- **Value Iteration**

Finds the optimal policy by iteratively estimating the utility

**INTERMISSION**

# **REINFORCEMENT LEARNING**

as Dynamic Programming



# REINFORCEMENT LEARNING RECAP



# REINFORCEMENT LEARNING RECAP

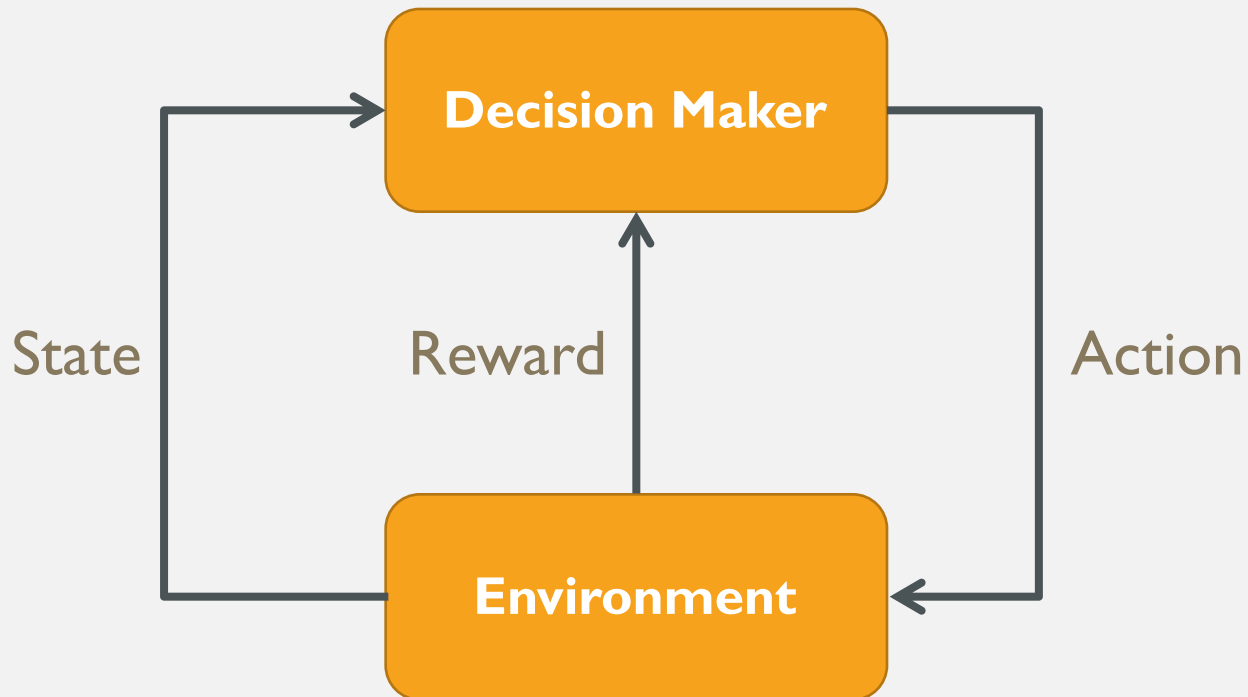


**Model-Based**

**Model-Free**

# REINFORCEMENT LEARNING RECAP

## Model-Based (MDP)



- Set of States:  $\mathcal{S} = \{s\}$
- Set of Actions:  $\mathcal{A} = \{a\}, a: \mathcal{S} \rightarrow \mathcal{S}$
- Reward Function:  $R(s)$
- Transition Model:  
 $T(s, a, s')$  or  $P(s' | s, a)$
- Policy:  $\pi(s) \rightarrow a$

# REINFORCEMENT LEARNING RECAP

## Model-Free



- Set of States:  $\mathcal{S} = \{s\}$
- Set of Actions:  $\mathcal{A} = \{a\}, a: \mathcal{S} \rightarrow \mathcal{S}$

There is no explicit model for  $R$  and  $T$   
We instead have explicit transition data:

$$\langle s, a, s', r \rangle$$

- Policy:  $\pi(s) \rightarrow a$

# MODEL-FREE REINFORCEMENT LEARNING

- Temporal Difference Learning
- Q-Learning

# TEMPORAL DIFFERENCE (TD) LEARNING

$$U(\mathbf{s}) = R(\mathbf{s}, \mathbf{a}) + \gamma \sum_{\mathbf{s}'} T(\mathbf{s}, \mathbf{a}, \mathbf{s}') U(\mathbf{s}')$$

Upon an action  $\mathbf{a} = \pi(\mathbf{s})$

For all  $\mathbf{s}'$ , successor of  $\mathbf{s}$ ,  $U(\mathbf{s})$  must be “in between”

a) the new value considering only  $\mathbf{s}'$ :  $R(\mathbf{s}) + \gamma U(\mathbf{s}')$

b) the old value  $U(\mathbf{s})$

# TEMPORAL DIFFERENCE (TD) LEARNING

$$U(\mathbf{s}) = R(\mathbf{s}, \mathbf{a}) + \gamma \sum_{\mathbf{s}'} T(\mathbf{s}, \mathbf{a}, \mathbf{s}') U(\mathbf{s}')$$

Upon an action  $\mathbf{a} = \pi(\mathbf{s})$

For all  $\mathbf{s}'$ , successor of  $\mathbf{s}$ ,  $U(\mathbf{s})$  must be “in between”

a) the new value considering only  $\mathbf{s}'$ :  $R(\mathbf{s}) + \gamma U(\mathbf{s}')$

b) the old value  $U(\mathbf{s})$



The Notion of Temporal Difference

# TEMPORAL DIFFERENCE (TD) LEARNING

$$U(\mathbf{s}) := (1 - \alpha) U(\mathbf{s}) + \alpha (R(\mathbf{s}) + \gamma U(\mathbf{s}'))$$

The new approximation of  $U(\mathbf{s})$  using  $\alpha$  when moving from a state  $\mathbf{s}$  to another state  $\mathbf{s}'$

Rearrange the above equation to update  $U(\mathbf{s})$ :

$$U(\mathbf{s}) := U(\mathbf{s}) + \alpha (R(\mathbf{s}) + \gamma U(\mathbf{s}') - U(\mathbf{s}))$$

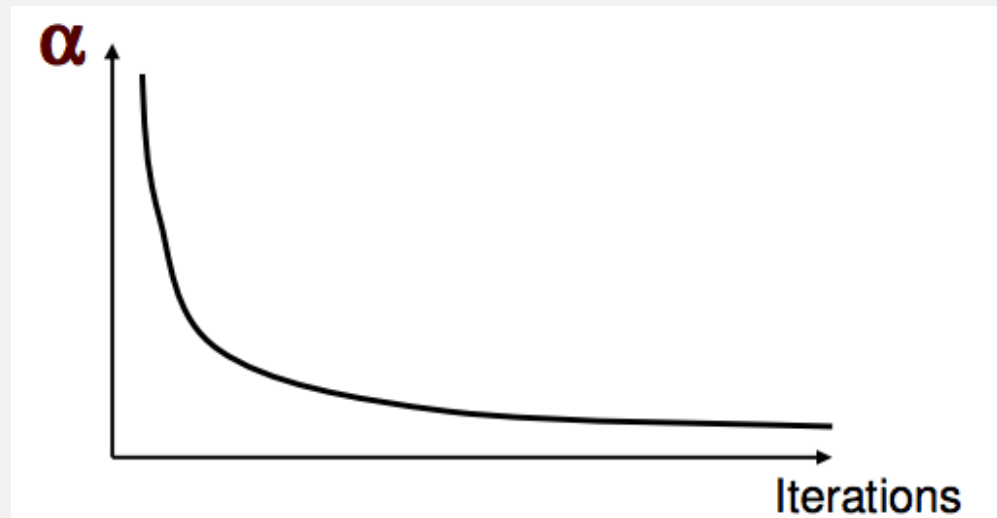


# TEMPORAL DIFFERENCE (TD) LEARNING

$$U(\mathbf{s}) := U(\mathbf{s}) + \alpha (R(\mathbf{s}) + \gamma U(\mathbf{s}') - U(\mathbf{s}))$$



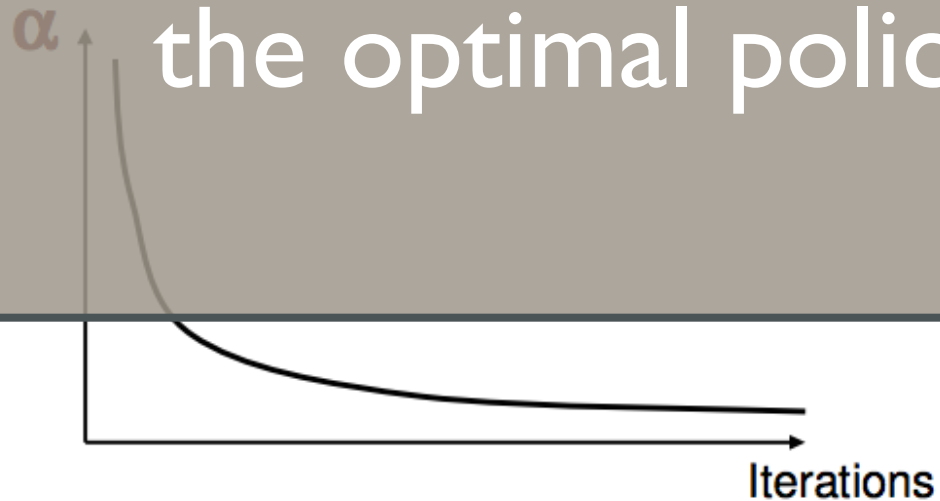
Learning Rate:  $0 \leq \alpha < 1$



# TEMPORAL DIFFERENCE (TD) LEARNING

$$U(s) := U(s) + \alpha (R(s) + \gamma U(s') - U(s))$$

But this says nothing about  
the optimal policy,  $\pi^*$



# Q-LEARNING

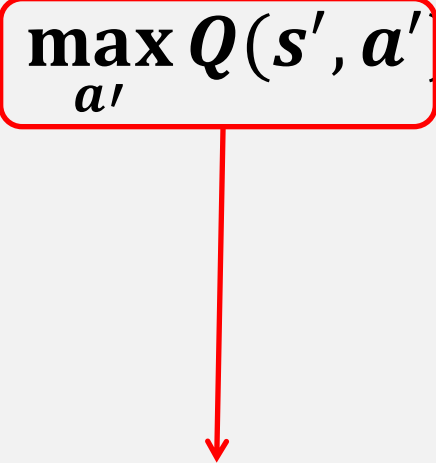
$$U(\mathbf{s}) := U(\mathbf{s}) + \alpha (R(\mathbf{s}) + \gamma U(\mathbf{s}') - U(\mathbf{s}))$$

Policy is about taking an action  $\mathbf{a}$  from a state  $\mathbf{s}$   
--> Let's introduce action to  $U(\mathbf{s})!$

$Q(\mathbf{s}, \mathbf{a})$ : Expected sum of future discounted rewards  
after taking action  $\mathbf{a}$  at state  $\mathbf{s}$

$$Q(\mathbf{s}, \mathbf{a}) := Q(\mathbf{s}, \mathbf{a}) + \alpha (R(\mathbf{s}) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}') - Q(\mathbf{s}, \mathbf{a}))$$

# Q-LEARNING

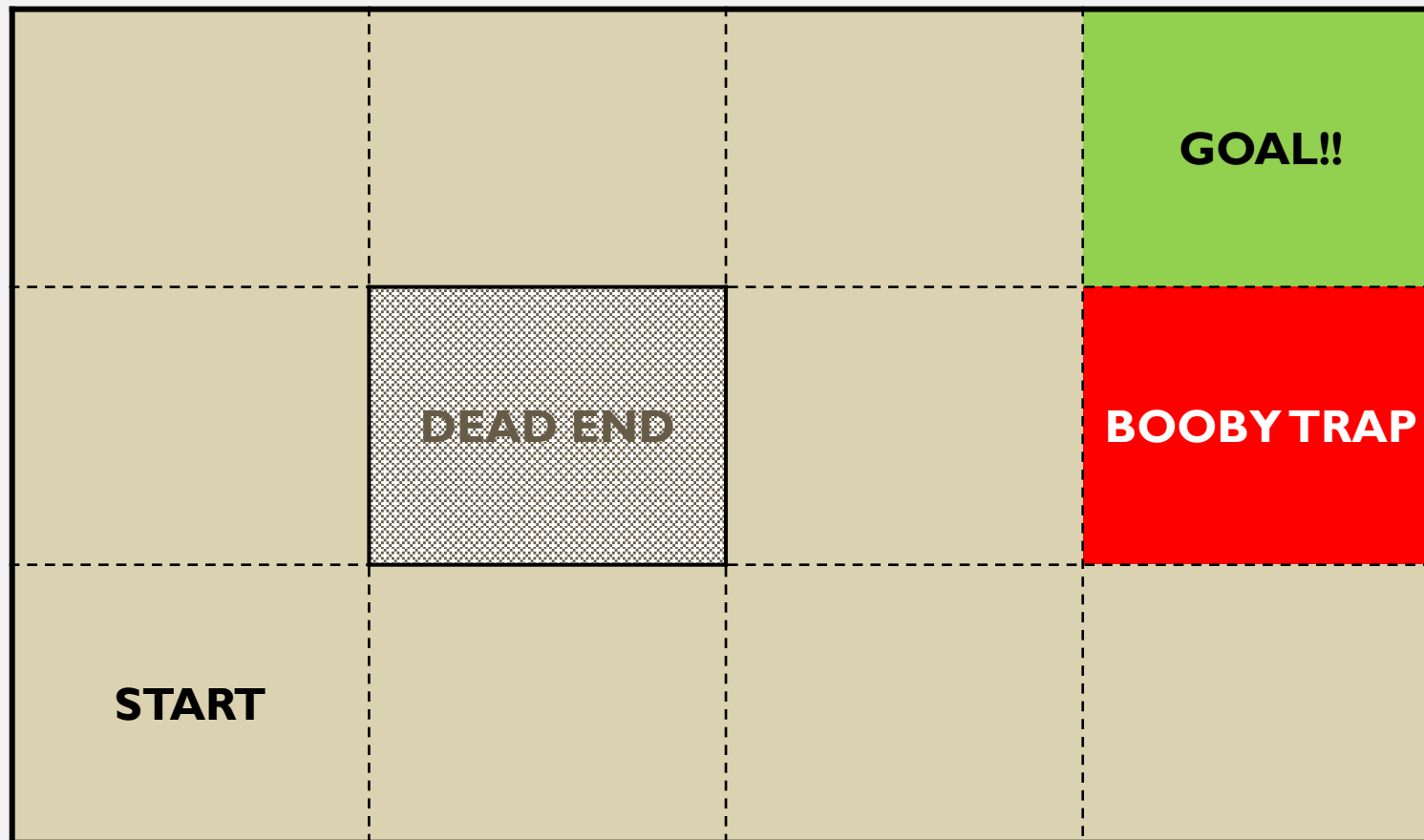
$$Q(s, a) := Q(s, a) + \alpha(R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$


Finding the optimal policy is now a matter of finding the actions for the states in maximizing the Q-value:

$$\pi(s) = \operatorname{argmax}_a Q(s, a)$$

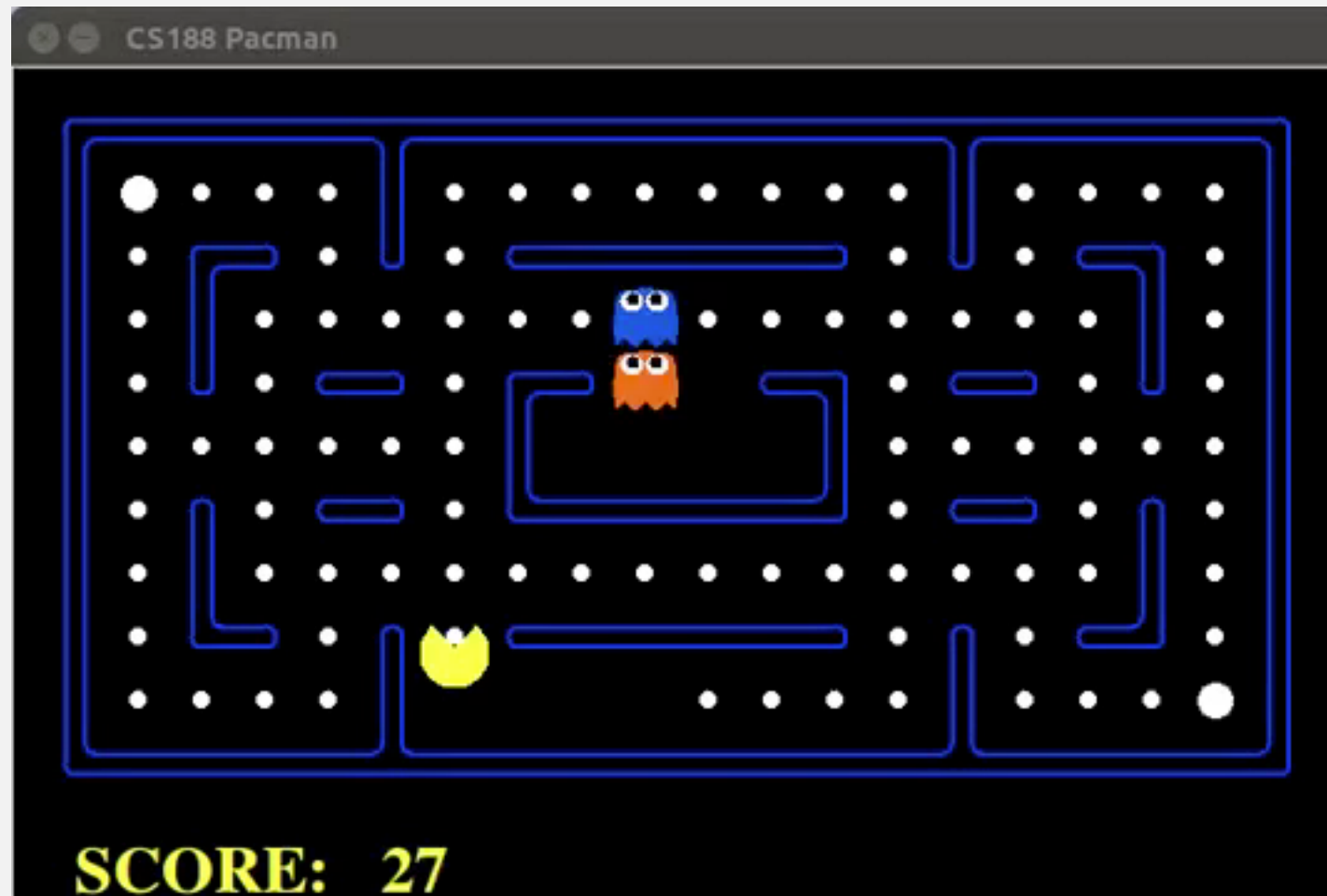
**EXAMPLE**

# SCENARIO



You can move UP, DOWN, LEFT, RIGHT

# PACMAN

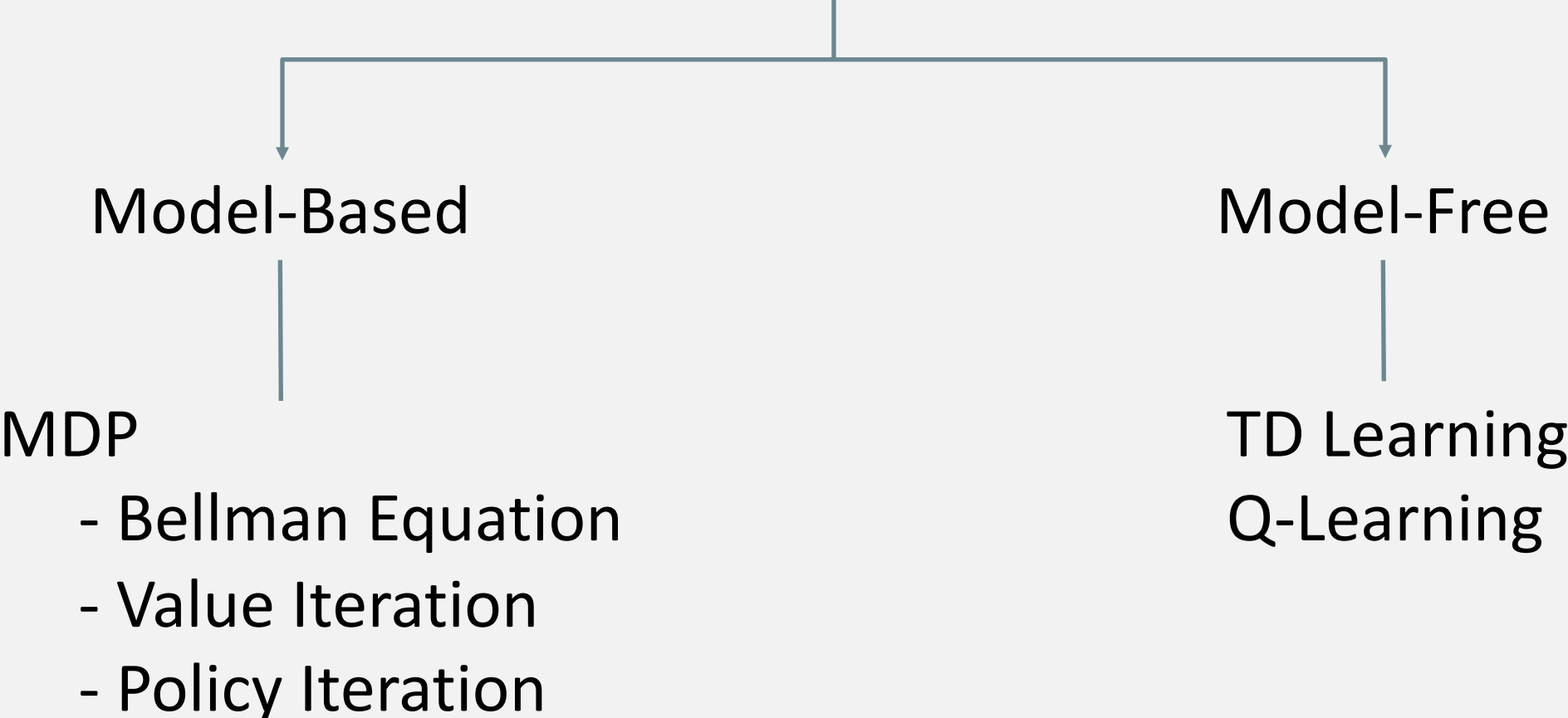


# **SUMMARY**



# SUMMARY

## Reinforcement Learning



# **DISCUSSION**

# QUESTIONS FOR THOUGHT

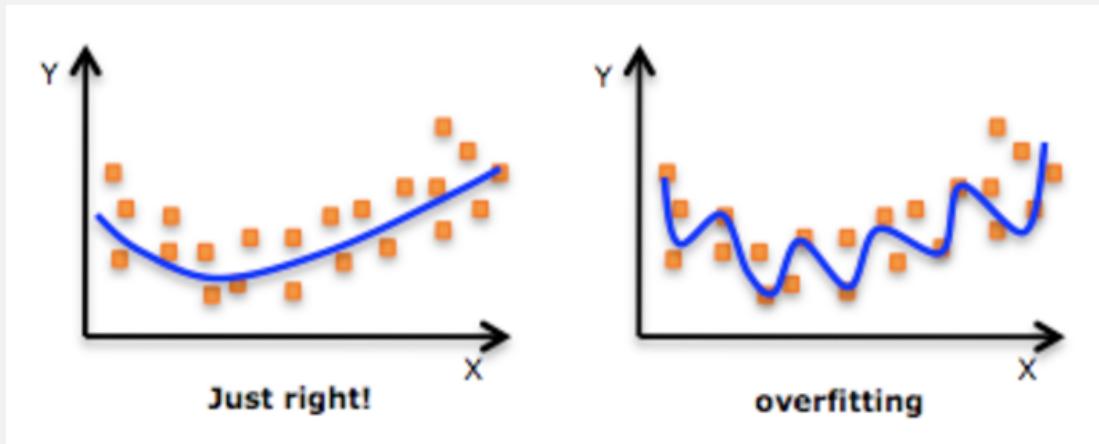
**Q1. What if the problem we are trying to solve is large scale?**

⇒ State-space is large!

⇒ ex. value iteration may take forever!

⇒ examining the entire training dataset is endless!

**Q2. What if our solution “overfits” data given by the environment?**



# **APPROXIMATE DYNAMIC PROGRAMMING**

# LINEAR APPROXIMATION APPROACH

**The curse-of-dimensionality problem:** the exponential growth of computational complexity with increasing dimensionality of the state space.

- a) **Abandon the idealized notion of optimality**  
and be content with a **suboptimal** solution 😊
- b) **Approximate functions for data:  $\langle s, a, s', r \rangle$**   
- can prevent overfitting

**Q & A**

# **APPENDIX**

# NOTATION & TERMINOLOGY

- We found that the notations used in the book are tricky to understand, and concluded to use the equivalent alternatives that were easier.

## <TEXTBOOK>

- $i \sim X_n$ : State
- $\mu(i)$ : Policy
- $J^\pi(i)$ : Cost-to-go function
- $g(X_n, \mu_n, X_{n+1})$ : Transition Cost
- $p_{ij}(\mu(i))$ : Transition Probability

## <THE SLIDES>

- $s \in S$ : State
- $\pi(s)$ : Policy
- $U(s)$ : Utility
- $R(s, a)$ : Reward (or  $R(s)$ )
- $T(s, a, s')$ : Transition Probability